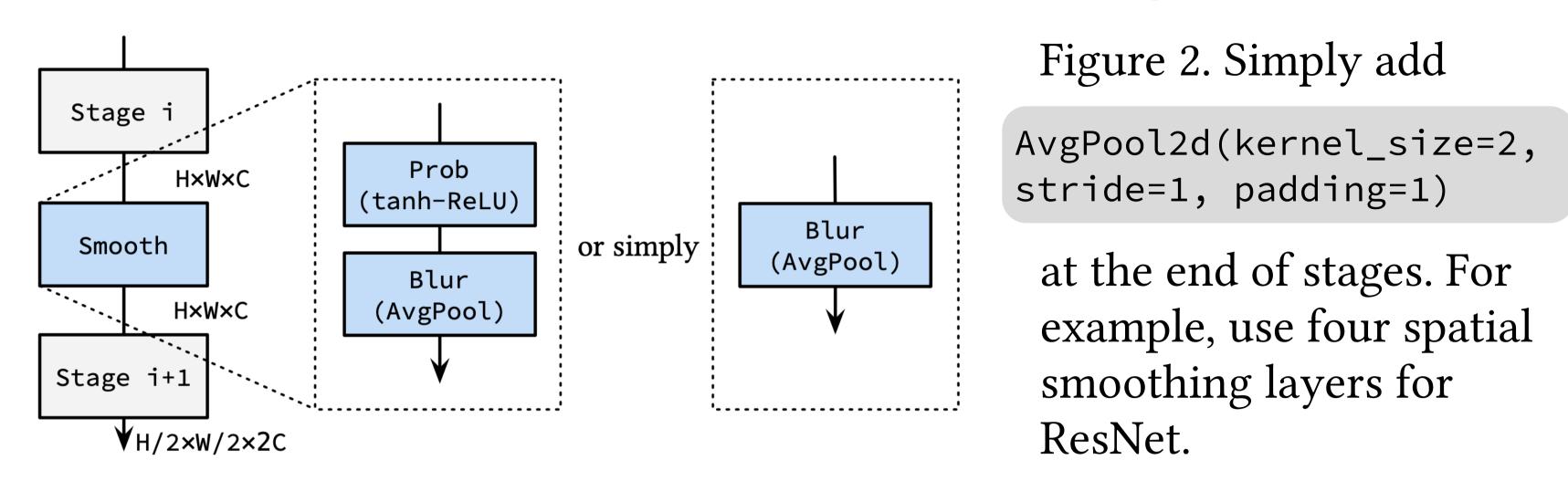
Blurs Behave Like Ensembles

We introduce a novel ensemble method, "spatial ensemble". Spatial ensemble is an extremely easy-to-implement method and improve accuracy, uncertainty, and robustness without increasing inference time.

I. What Is a Spatial Ensemble?

Figure 1. Spatial ensemble, or spatial smoothing, is a method that aggregate nearby feature maps. Please refer to the figure on the right (\searrow) .

II. How Can We Apply Spatial Smoothing to NNs?



III. Spatial Smoothing Improves MC Dropout

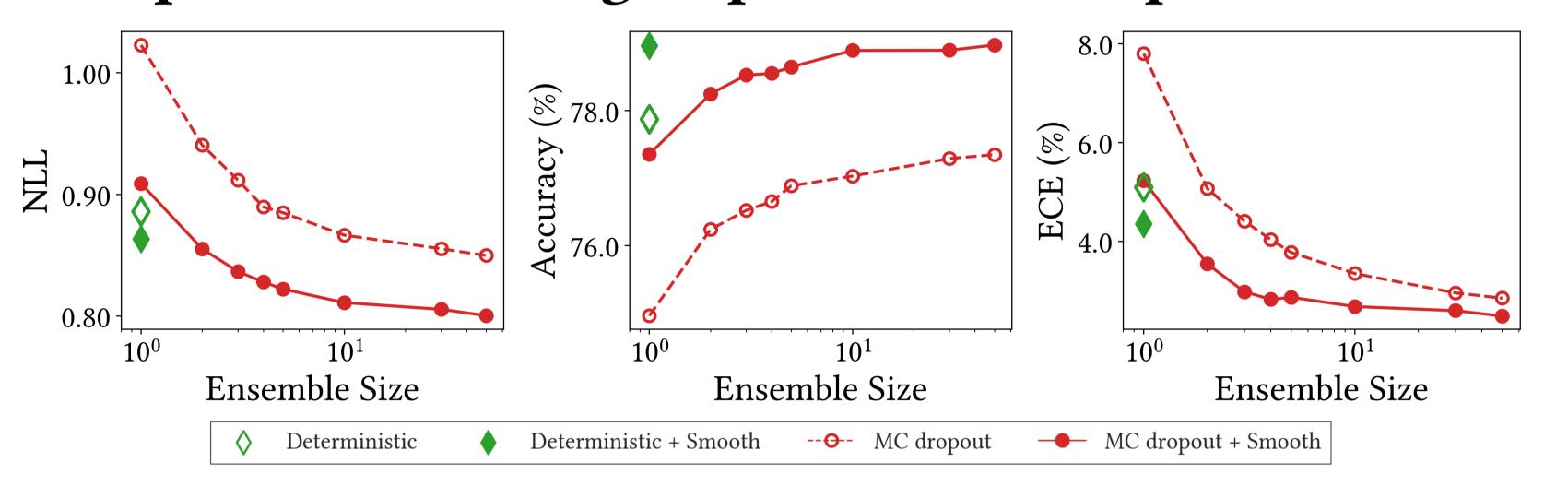


Figure 3. "MC dropout + spatial smoothing" is **25**× **faster** than canonical MC dropout with similar predictive performance. Moreover, spatial smoothing also can be applied to canonical deterministic NNs to improve the performances.

IV. How Does Spatial Smoothing Improve NNs?

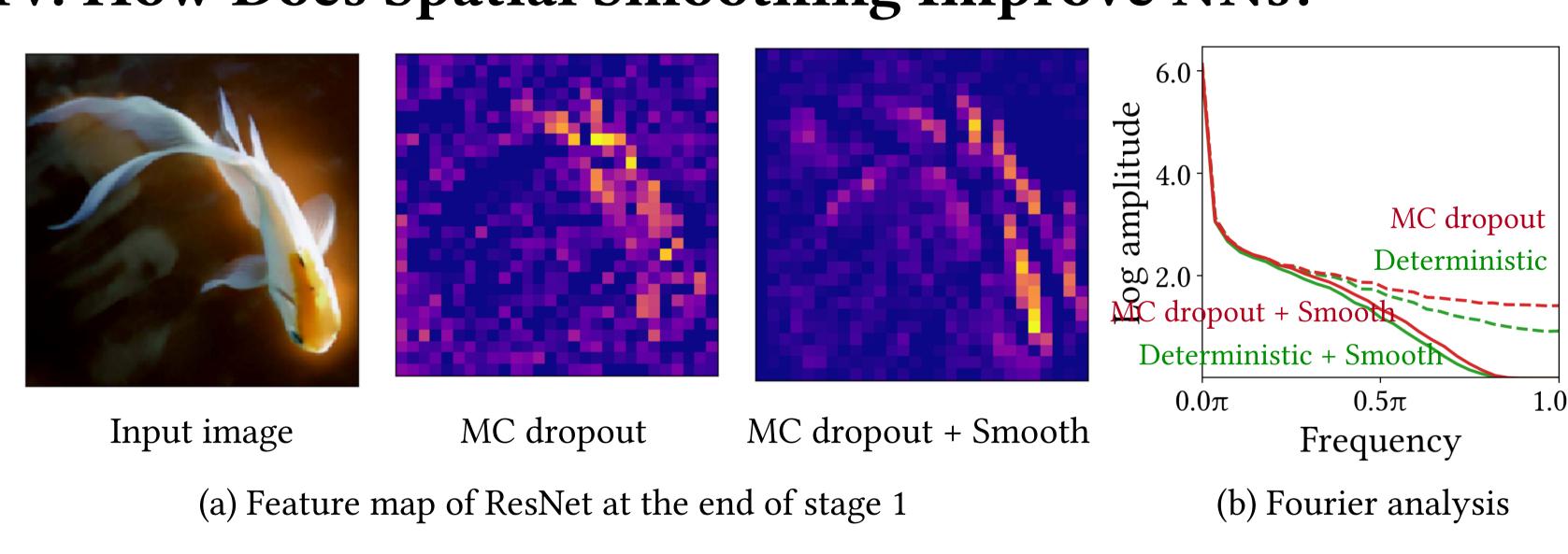


Figure 4. MC dropout adds high-frequency noises, and spatial smoothing **filters high-frequency signals** and stabilizes (denoises) feature maps.

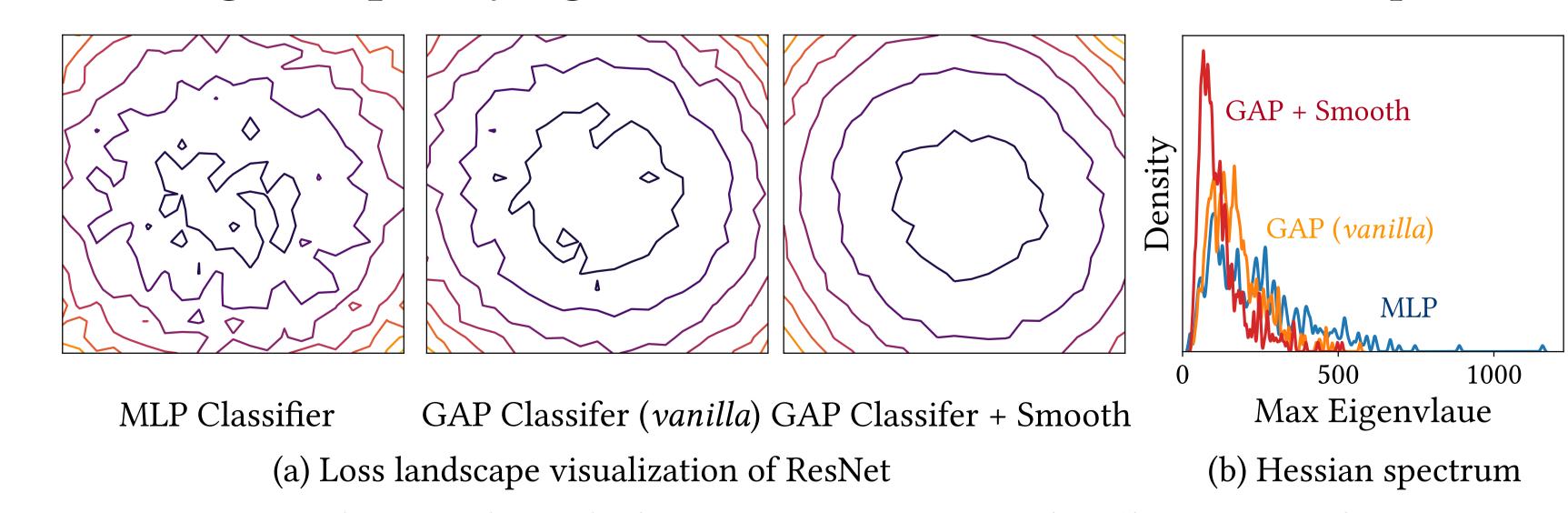


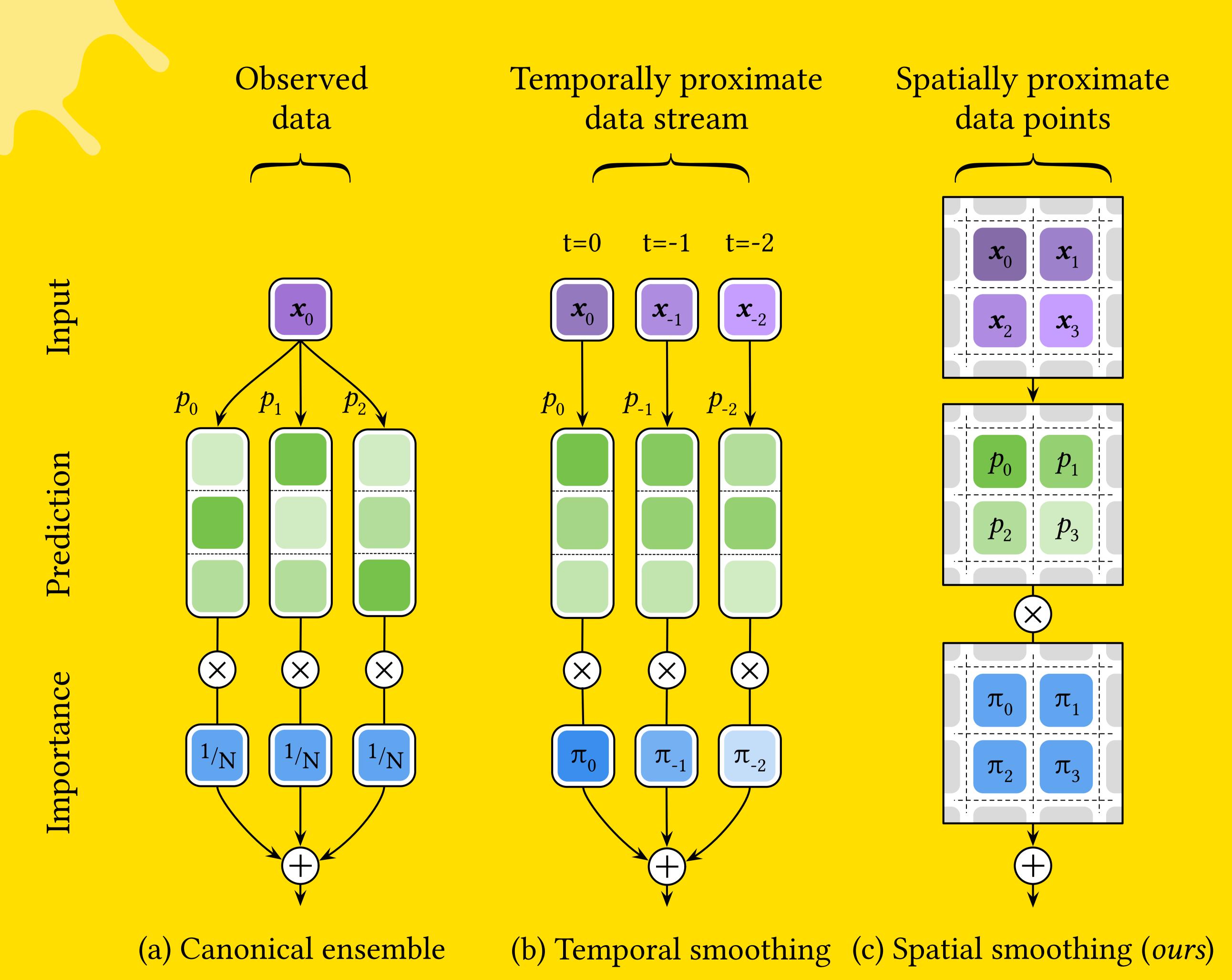
Figure 5. Spatial smoothing helps NN optimization by **flattening loss landscapes**. In addition, GAP is an extreme case of spatial smoothing.

V. How Can We Make Spatial Smoothing Trainable?

Self-attentions for computer vision, also known as Vision Transformers (ViTs), can be deemed as trainable importance-weighted ensembles of feature maps.



SPATIAL CONSISTENCY is important, and BLURS BEHAVE LIKE ENSEMBLES



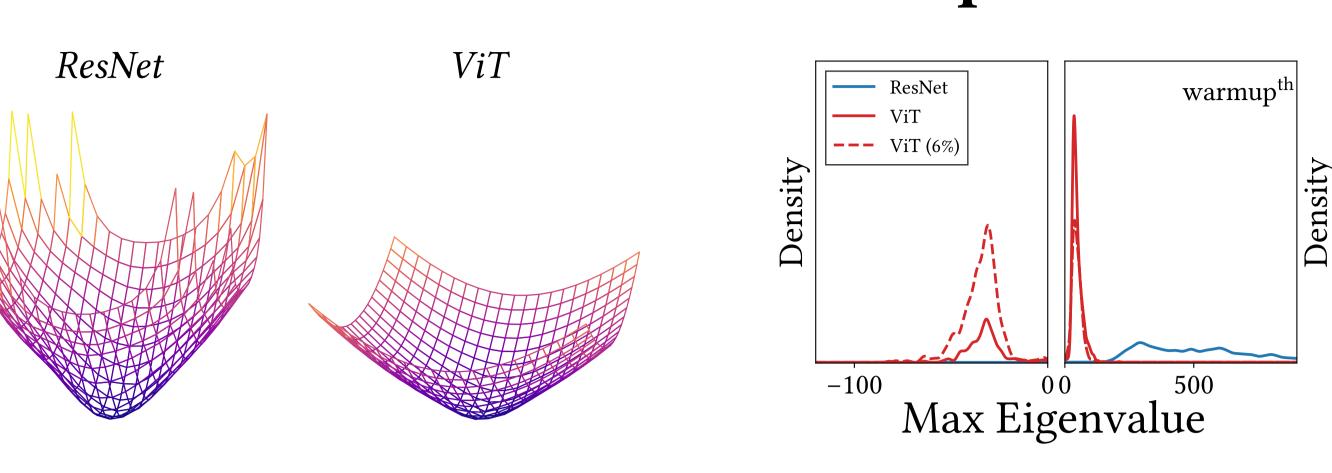
NAMUK PARK, SONGKUK KIM

NAVER AI Lab, Yonsei University

How Do Vision Transformers Work?

We provide an explanation of how *self-attentions* work by addressing them as *a trainable spatial smoothing* of feature maps, because the formulation suggests that self-attentions average feature map values with the positive importance-weights.

I. Self-Attentions Flatten Loss Landscapes



(a) Loss landscape visualization

(b) Hessian max eigenvalue spectrum

Figure 6. Loss landscape visualization (*Left*) and Hessian max eigenvalue spectrum (*Right*) consistently shows that **ViT has a flatter loss than ResNet**. It suggests that self-attentions flatten loss landscapes.

II. Self-Attentions Are Low-Pass Filters

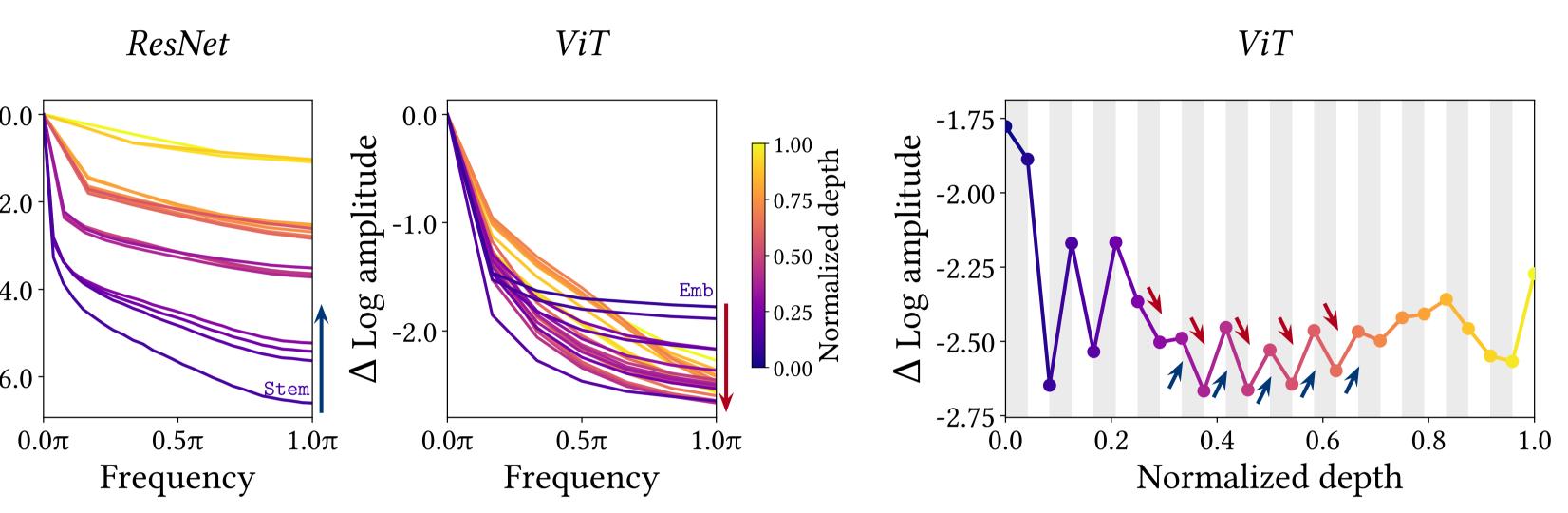
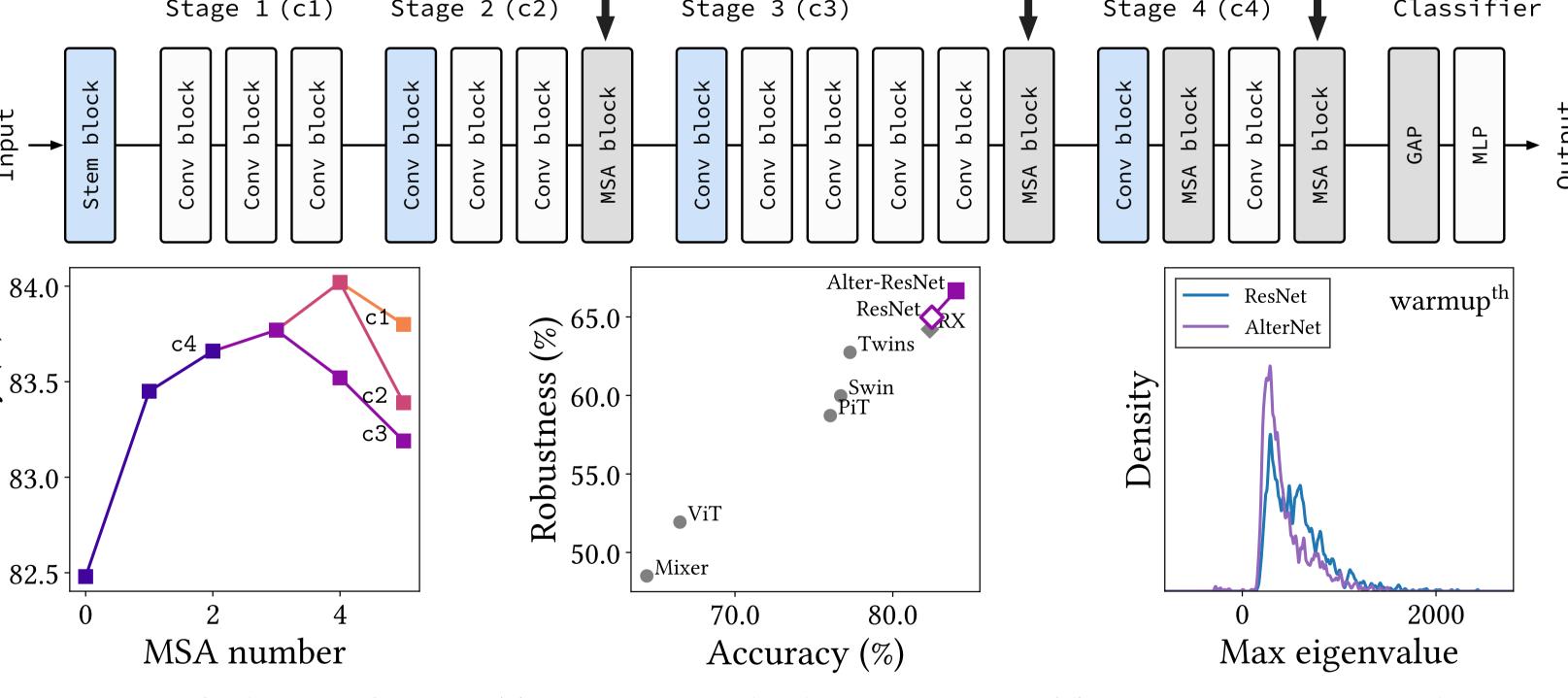


Figure 7. Relative log amplitudes of Fourier transformed feature map show that **ViT tend to reduces high-frequency signals, while ResNet amplify them**. *Right*: In ViT, MSAs (gray area) reduce the high-frequency (1.0π) component of feature map, and Conv/MLPs (white area) amplify it.

III. Self-Attentions at the End of a Stage Play a Key Role in Prediction



(b) Accuracy of AlterNet for MSA number

(c) Accuracy and robustness in a (d) Hessian max eigenvalue spectra small data regime (CIFAR-100)

Figure 8. Self-attentions closer to the end of a stage (not a model) and Convs at the beginning of a stage significantly improve the performance. AlterNet, a model in which Conv blocks at the end of a stage are replaced with MSA blocks, outperforms CNNs even in small data regimes.

In summary, Self-attentions have a spatial smoothing inductive bias that complements convolutions.

	Self-Attention	Convolution
Loss Landscape	Flat but non-convex	Convex but sharp
Fourier Analysis	Low-pass filter (shape-biased)	High-pass filter (texture-biased)
Best Practice	The end of a stage	The beginning of a stage

